# Sparse Matrices Beyond Solvers - Graphs, Biology, and Machine Learning (v2)

Aydın Buluç

Computational Research Division, LBNL

EECS Department, UC Berkeley

CS Summer Student Program
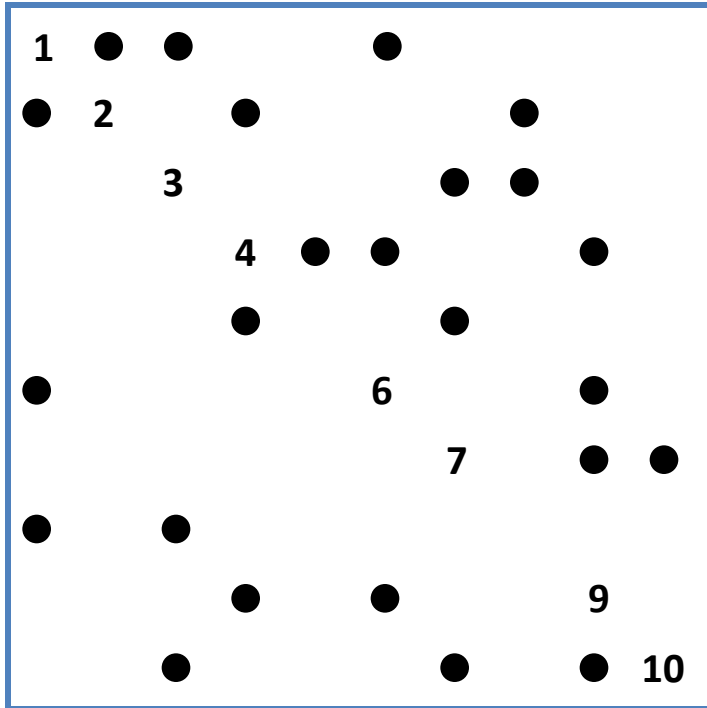
July 16, 2020

# Sparse Matrices

"I observed that most of the coefficients in our matrices were zero; i.e., the nonzeros were 'sparse' in the matrix, and that typically the triangular matrices associated with the forward and back solution provided by Gaussian elimination would remain sparse if pivot elements were chosen with care"
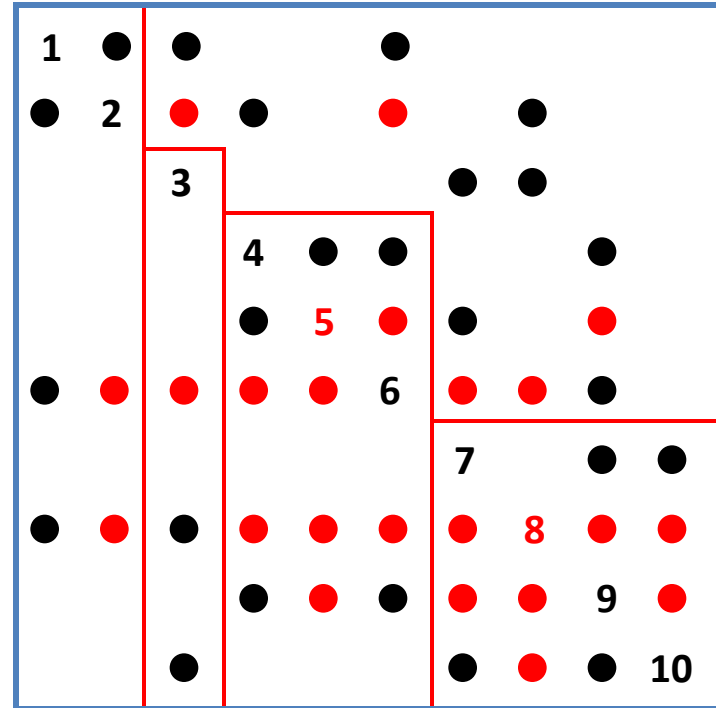
- Harry Markowitz, describing the 1950s work on portfolio theory that won the 1990 Nobel Prize for Economics
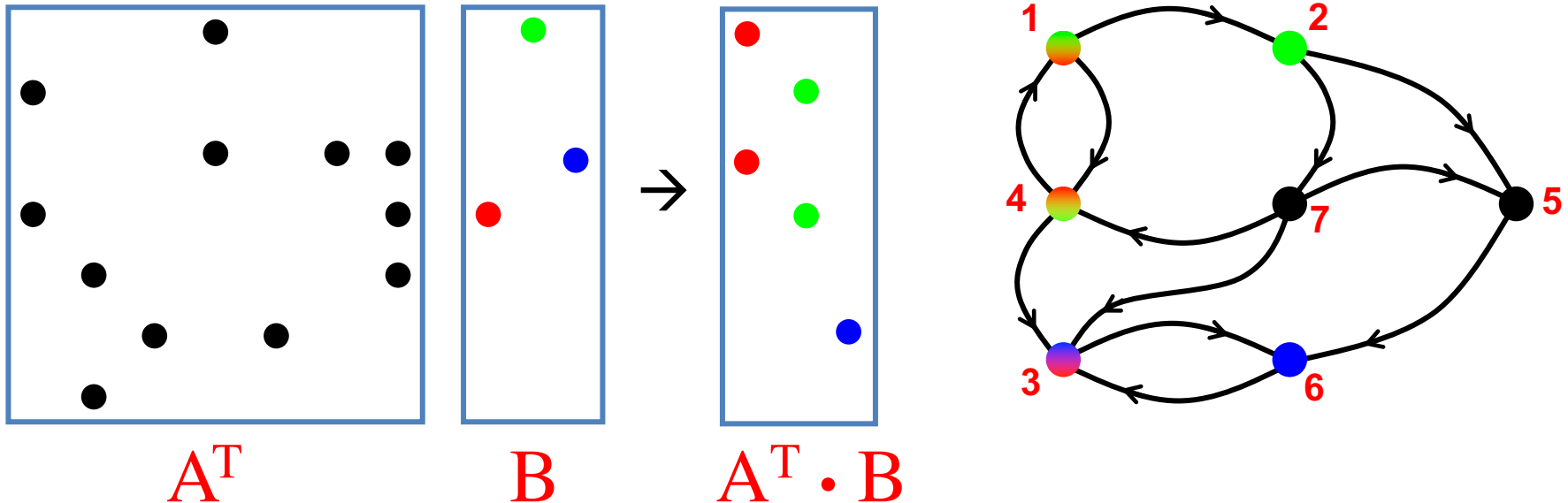
# Sparse Matrices



Original matrix A

Factors L+U

Original: Ax = b (hard to solve directly)

Factored: LUx = b (solvable by direct substitution)

# Graphs in the language of matrices
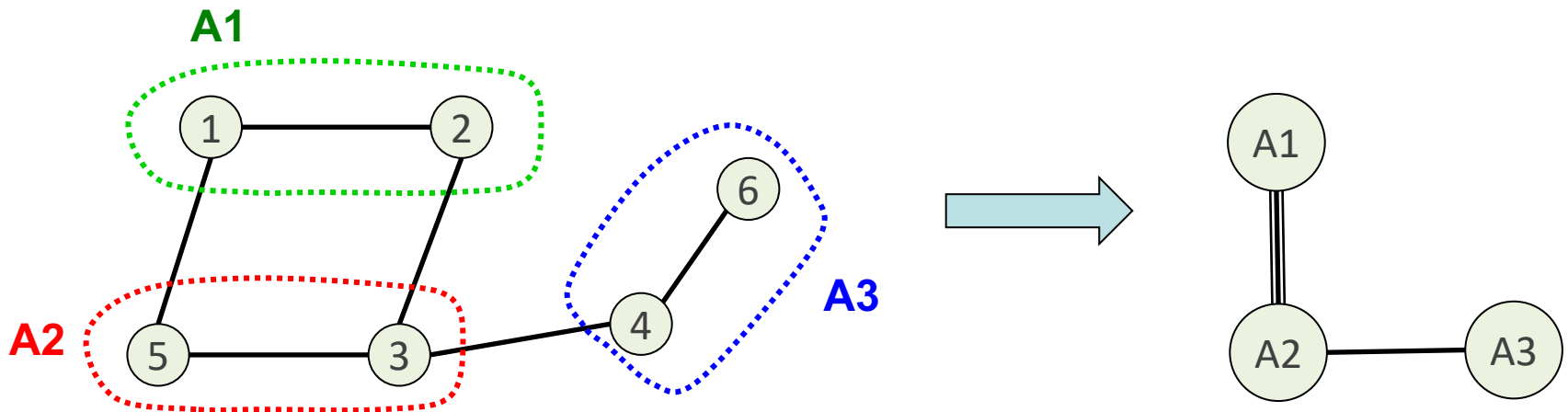


$$A^T \qquad B \qquad A^T \cdot B$$

- Sparse array representation => space efficient
- Sparse matrix-matrix multiplication => work efficient
- Three possible levels of parallelism:  searches, vertices, edges
- Highly-parallel implementation for Betweenness Centrality*

  *: A measure of influence in graphs, based on shortest paths

# Graph coarsening via sparse matrix-matrix products

Aydin Buluç and John R. Gilbert. Parallel sparse matrix-matrix multiplication and indexing: Implementation and experiments. *SIAM Journal of Scientific Computing (SISC), 2012*.
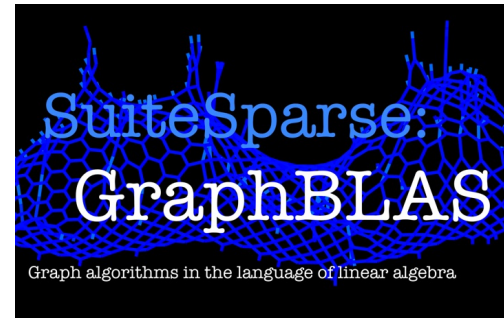
# The GraphBLAS effort

## Standards for Graph Algorithm Primitives

Tim Mattson (Intel Corporation), David Bader (Georgia Institute of Technology), Jon Berry (Sandia National Laboratory), Aydin Buluc (Lawrence Berkeley National Laboratory), Jack Dongarra (University of Tennessee), Christos Faloutsos (Carnegie Melon University), John Feo (Pacific Northwest National Laboratory), John Gilbert (University of California at Santa Barbara), Joseph Gonzalez (University of California at Berkeley), Bruce Hendrickson (Sandia National Laboratory), Jeremy Kepner (Massachusetts Institute of Technology), Charles Leiserson (Massachusetts Institute of Technology), Andrew Lumsdaine (Indiana University), David Padua (University of Illinois at Urbana-Champaign), Stephen Poole (Oak Ridge National Laboratory), Steve Reinhardt (Cray Corporation), Mike Stonebraker (Massachusetts Institute of Technology), Steve Wallach (Convey Corporation), Andrew Yoo (Lawrence Livermore National Laboratory)

*Abstract*-- It is our view that the state of the art in constructing a large collection of graph algorithms in terms of linear algebraic operations is mature enough to support the emergence of a standard set of primitive building blocks. This paper is a position paper defining the problem and announcing our intention to launch an open effort to define this standard.

- The GraphBLAS Forum: http://graphblas.org
- Graphs: Architectures, Programming, and Learning  (GrAPL @IPDPS): http://hpc.pnl.gov/grapl/

# SuiteSparse::GraphBLAS

- From Tim Davis (Texas A&M)
- First conforming implementation of C API
- Features [1]:
  - 960 semirings built in; also user-defined semirings
  - Fast incremental updates using non-blocking mode and "zombies"
  - Several sparse data structures & polyalgorithms under the hood
- Already multithreaded [2]
- Performance on graph benchmarks (e.g. triangles, k-truss) comparable to highly-tuned custom C code
- Included in Debian and Ubuntu Linux distributions
- Used as computational engine in commercial RedisGraph product

[1] Davis, Timothy A. "Algorithm 1000: SuiteSparse: GraphBLAS: Graph Algorithms in the Language of Sparse Linear Algebra." ACM Transactions on Mathematical Software (TOMS) 45.4 (2019): 44.
[2] Aznaveh, Mohsen, et al. "Parallel GraphBLAS with OpenMP." CSC20, SIAM Workshop on Combinatorial Scientific Computing. SIAM. 2020.

# GraphBLAS C API Spec ([http://graphblas.org](http://graphblas.org))

- **Goal:** A crucial piece of the GraphBLAS effort is to translate the mathematical specification to an actual Application Programming Interface (API) that
  - i. is faithful to the mathematics as much as possible, and
  - ii. enables efficient implementations on modern hardware.
- **Impact:** All graph and machine learning algorithms that can be expressed in the language of linear algebra
- **Innovation:** Function signatures (e.g. mxm, vxm, assign, extract), parallelism constructs (blocking v. non-blocking), fundamental objects (masks, matrices, vectors, descriptors), a hierarchy of algebras (functions, monoids, and semiring)

```
GrB_info GrB_mxm(GrB_Matrix          *C,       // destination
         const GrB_Matrix             Mask,
         const GrB_BinaryOp           accum,
         const GrB_Semiring           op,
         const GrB_Matrix             A,
         const GrB_Matrix             B
     [, const Descriptor             desc]);
```
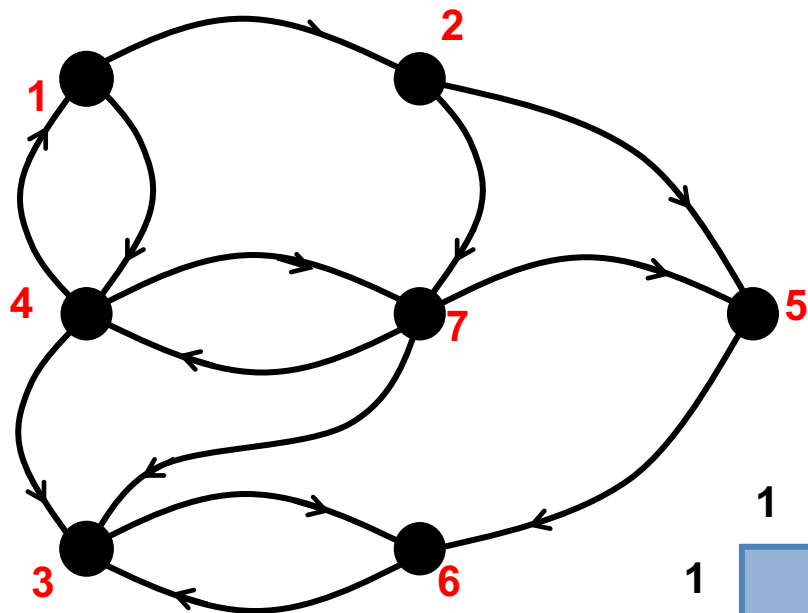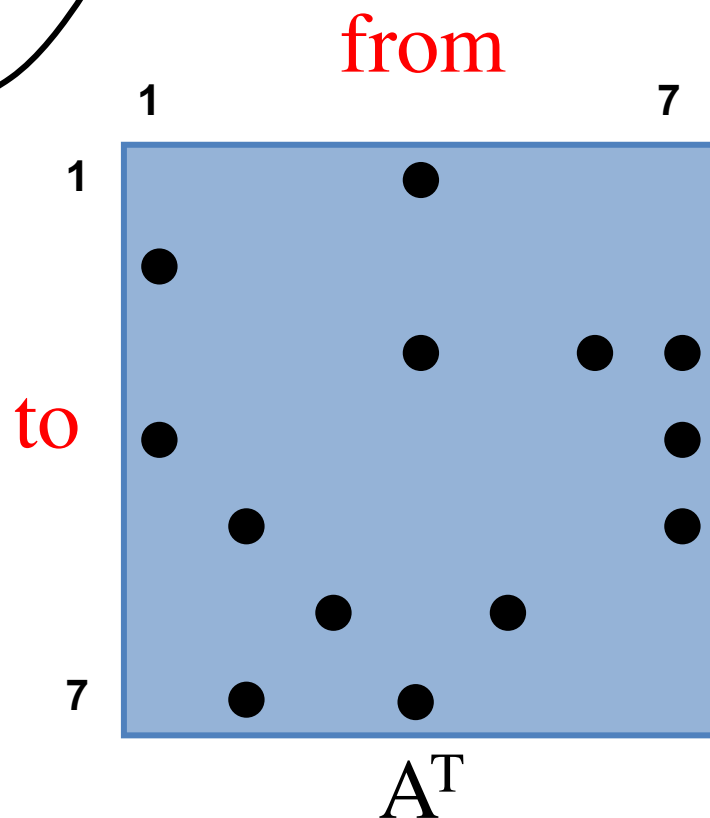
$$C(\neg M) \oplus= A^T \oplus.\otimes B^T$$

A.Buluç, T. Mattson, S. McMillan, J. Moreira, C. Yang. "The GraphBLAS C API Specification", version 1.3.0

# Examples of semirings in graph algorithms

| | |
|---|---|
| Real field: **(R, +, x)** | Classical numerical linear algebra |
| Boolean algebra: **({0 1}, \|, &)** | Graph connectivity |
| Tropical semiring: **(R U {∞}, min, +)** | Shortest paths |
| **(S, select, select)** | Select subgraph, or contract nodes to form quotient graph |
| (edge/vertex attributes, vertex data aggregation, edge data processing) | Schema for user-specified computation at vertices and edges |
| **(R, max, +)** | Graph matching &network alignment |
| **(R, min, times)** | Maximal independent set |

- **Shortened semiring notation: (Set, Add, Multiply)**. Both identities omitted.
- **Add:** Traverses edges, **Multiply:** Combines edges/paths at a vertex
- Neither add nor multiply needs to have an inverse.
- Both **add** and **multiply** are **associative**, **multiply distributes** over **add**

Breadth-first search in the language of matrices

Particular semiring operations:
**Multiply:** select2nd
**Add:** minimum

from

to

parents:

$A^T$   $X$   $A^TX$

# Input sparsity

- What was the cost of that $A^T x$ in the previous slide?
- If x is dense, it is O(nnz(A)) = O(m) where m=#edges
- **If x is sparse**, it is

$$\sum_{i:x_i \neq 0} nnz(A_{i:})$$

- ***Over all iterations*** of BFS, the cost sums up to O(nnz(A)), because no $x_i$ appears twice in the input.
- Note that this is ***optimal*** for conventional (top-down) BFS
- Many people outside the community miss this observation and ***mistakenly think*** SpMV based BFS is suboptimal by a factor of the graph diameter.

Select vertex with
__minimum__ label as parent

from

to

parents:

$A^T$

$X$

$A^T X$

- Masks avoid formation of temporaries and can enable automatic direction optimization
- These footballs are nonzeros that are **masked out** by the parents array

parents:

from

to

$A^T$     $X$     $A^TX$

1     7

from

1      7

to

$A^T$     X     $A^TX$

# GraphBLAST

- First "high-performance" GraphBLAS implementation on the GPU
- Optimized to take advantage of both input and output sparsity
- Automatic direction-optimization through the use of masks
- Competitive with fastest GPU (Gunrock) and CPU (Ligra) codes
- Outperforms multithreaded SuiteSparse::GraphBLAS

Design principles:
1. Exploit input sparsity => direction-optimization
2. Exploit output sparsity => masking
3. Proper load-balancing => key for GPU implementations

Extensively evaluated on (more implemented, google for github repo)
- Breadth-first-search (BFS)
- Single-source shortest-path (SSSP)
- PageRank (PR)
- Triangle counting (TC)        https://github.com/gunrock/graphblast

Yang, B., Owens, "GraphBLAST: A High-Performance Linear Algebra-based Graph Framework on the GPU", arXiv

# Kernel methods in Machine Learning

| A **kernel** is a function that | Implicitly transforms raw data into high-dimensional feature vectors via a **feature map**; and then | Returns an **inner product** between the feature vectors. Must be **positive-definite.** |
| --- | --- | --- |
| A **kernel** is useful for | **Factor out** knowledge on data representation from downstream algorithms, | Exploit **infinite dimensionality and nonlinear** feature spaces. |
| **Kernels** are used in | Support vector machine (SVM), Gaussian process regression (GPR), Kernel principal component analysis (kPCA), etc. | |



(a)  (b)

Figure source: Russell & Norvig

The circular decision boundary in 2D (a) becomes a linear boundary in 3D (b) using the following transformation: $\phi(x_1, x_2) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$

# Marginalized Graph Kernels

The inner product between two graphs is the statistical average of the inner product of simultaneous random walk paths on the two graphs.

Graph A

Graph A
0.9
0.4    0.6
0.9

Use edge weight to set transition probability

Sample paths

Length=1
$p = 0.4$
$p = 0.6$

Length=2
$p = 0.4 \times 0.9 = 0.36$
$p = 0.6 \times 0.9 = 0.54$

Compare

Graph B

Graph B
0.5
0.2    0.3    0.3
0.4   0.7   0.2   0.6
0.5

$p = 0.2$
$p = 0.3$
$p = 0.5$

$p = 0.2 \times 0.5 = 0.10$
$p = 0.2 \times 0.4 = 0.08$
$p = 0.3 \times 0.3 = 0.09$
$p = 0.3 \times 0.6 = 0.18$
$p = 0.5 \times 0.7 = 0.35$
$p = 0.5 \times 0.6 = 0.30$

$$K(G, G') = \mathbf{p}_{\times}^{\mathsf{T}} \left( \mathbf{D}_{\times} \mathbf{V}_{\times}^{-1} - \mathbf{A}_{\times} \odot \mathbf{E}_{\times} \right)^{-1} \mathbf{D}_{\times} \mathbf{q}_{\times}$$

degree    vertex label    adjacency    edge label

SPD system to solve

The marginalized graph kernel in linear algebra form represents a modified graph Laplacian

**Streaming Kronecker matrix-vector multiplication**
- Regenerates the product linear system on the fly by streaming 8-by-8 tiles.
- Tiles staged in shared memory.
- Trade FLOPS for GB/s, but asymptotic arithmetic complexity stays the same.



$$
\begin{aligned}
&1 \quad \text{function } \mathrm{CG4GK}(\mathbf{d},\mathbf{d}',\mathbf{v},\mathbf{v}',\mathbf{A},\mathbf{A}',\mathbf{E},\mathbf{E}',\ \mathbf{q},\mathbf{q}')\\
&2 \qquad \mathbf{M} \leftarrow \mathbf{diag}\left[(\mathbf{d}\otimes\mathbf{d}')\odot(\mathbf{v}\overset{\kappa}{\otimes}\mathbf{v}')^{-1}\right] \qquad |{+}|\\
&3 \qquad \mathbf{x} \leftarrow \mathbf{0} \qquad |{+}|\\
&4 \qquad \mathbf{r} \leftarrow (\mathbf{d}\otimes\mathbf{d}')\cdot(\mathbf{q}\otimes\mathbf{q}') \qquad \boxtimes\cdot|\\
&5 \qquad \mathbf{z} \leftarrow \mathbf{v}\overset{\kappa}{\otimes}\mathbf{v}' \qquad |{+}|\\
&6 \qquad \mathbf{p} \leftarrow \mathbf{z} \qquad |{+}|\\
&7 \qquad \rho \leftarrow \mathbf{r}^{\mathsf{T}}\mathbf{z} \qquad |^{\mathsf{T}}\cdot|\\
&8 \qquad \text{repeat}\\
&9 \qquad\qquad \mathbf{a} \leftarrow (\mathbf{d}\otimes\mathbf{d}')\odot(\mathbf{v}\overset{\kappa}{\otimes}\mathbf{v}')^{-1}\cdot\mathbf{p} \qquad \boxtimes\cdot|\\
&10 \qquad\qquad\quad +(\mathbf{A}\otimes\mathbf{A}')\odot(\mathbf{E}\overset{\kappa}{\otimes}\mathbf{E}')\cdot\mathbf{p} \qquad \boxtimes\cdot|\\
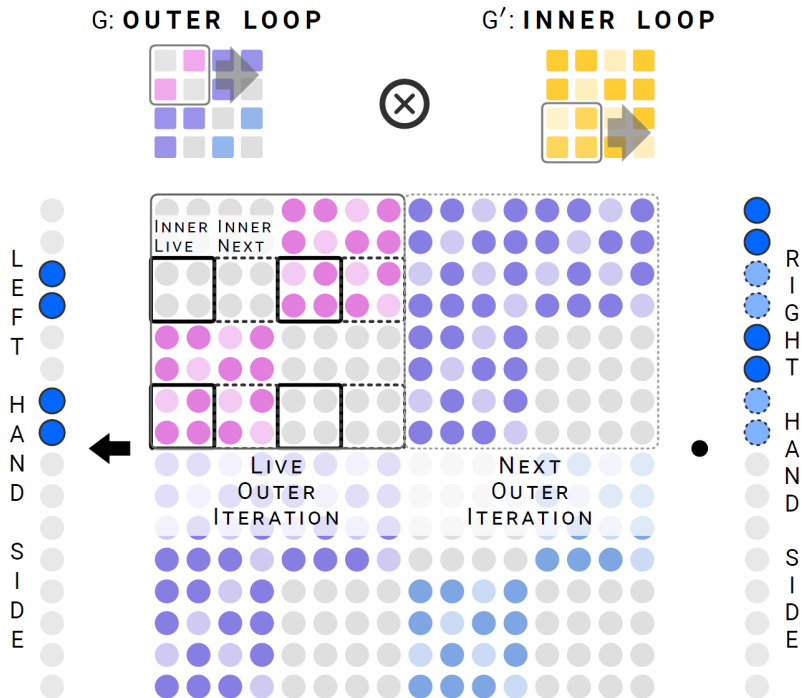&11 \qquad\qquad \alpha \leftarrow \rho/(\mathbf{p}^{\mathsf{T}}\mathbf{a}) \qquad |^{\mathsf{T}}\cdot|\\
&12 \qquad\qquad \mathbf{x} \leftarrow \mathbf{x}+\alpha\mathbf{p} \qquad |{+}|\\
&13 \qquad\qquad \mathbf{r} \leftarrow \mathbf{r}-\alpha\mathbf{a} \qquad |{+}|\\
&14 \qquad\qquad \mathbf{z} \leftarrow \mathbf{M}^{-1}\mathbf{r} \qquad |{+}|\\
&15 \qquad\qquad \rho' \leftarrow \mathbf{r}^{\mathsf{T}}\mathbf{z} \qquad |^{\mathsf{T}}\cdot|\\
&16 \qquad\qquad \beta \leftarrow \rho'/\rho\\
&17 \qquad\qquad \mathbf{p} \leftarrow \mathbf{z}+\beta\mathbf{p} \qquad |{+}|\\
&18 \qquad\qquad \rho \leftarrow \rho'\\
&19 \qquad \text{until } \mathbf{r}^{\mathsf{T}}\mathbf{r} < \epsilon\\
&20 \qquad \text{return } \mathbf{x}
\end{aligned}
$$

# Exploiting Sparsity

- Most discrete systems have natural sparsity (e.g. not all atoms are connected).
- 2-level sparsity exploitation:
    i.   Outer level: retain only non-empty tiles
    ii.  Inner level: use bitmap + compact storage format
- Packing into compact format: on CPU as a preprocessing step
- Unpacking for Streaming Kronxv: on GPU using bit magic + warp intrinsics
- Partition-based graph ordering reduces # non-empty tiles
    ☛ Cost easily amortized because we reorder each graph, not their product

EMPTY TILE DISCARDED     NON-EMPTY TILE COMPRESSED     DENSE STORAGE

| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

64-bit integer **nzmask**

0b1000001000001000100100100011101010010010011001100000011000000

0x0303324AE4122041

BITMAP

| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |

Non-empty tile compressed grid:

|   |   |   | K |   |   |
|   |   | H | L |   | R |
|   | E |   | M | Q |   |
|   | F |   |   | O |   |
|   |   | I |   |   |   |
|   |   | N |   |   |   |
| A | C | G | J |   | P |
| B | D |   |   |   | S |

# Performance of the Graph Kernel

## Time-to-solution comparison with GraKeL and GraphKernels

■ Present   ■ GraKeL   ■ GraphKernels

**DrugBank**
- Present: 172 seconds
- GraKeL: 6461× — 12.9 days
- GraphKernels: 998× — 2.0 days

**PDB**
- Present: 153 seconds
- GraKeL: 3297× — 5.8 days
- GraphKernels: 12430× — 22.0 days

Axis: 1 second | 1 min | 1 hour | 1 day | 1 month

GraKeL: Cython, multi-threading
GraphKernels: Python, no parallelization

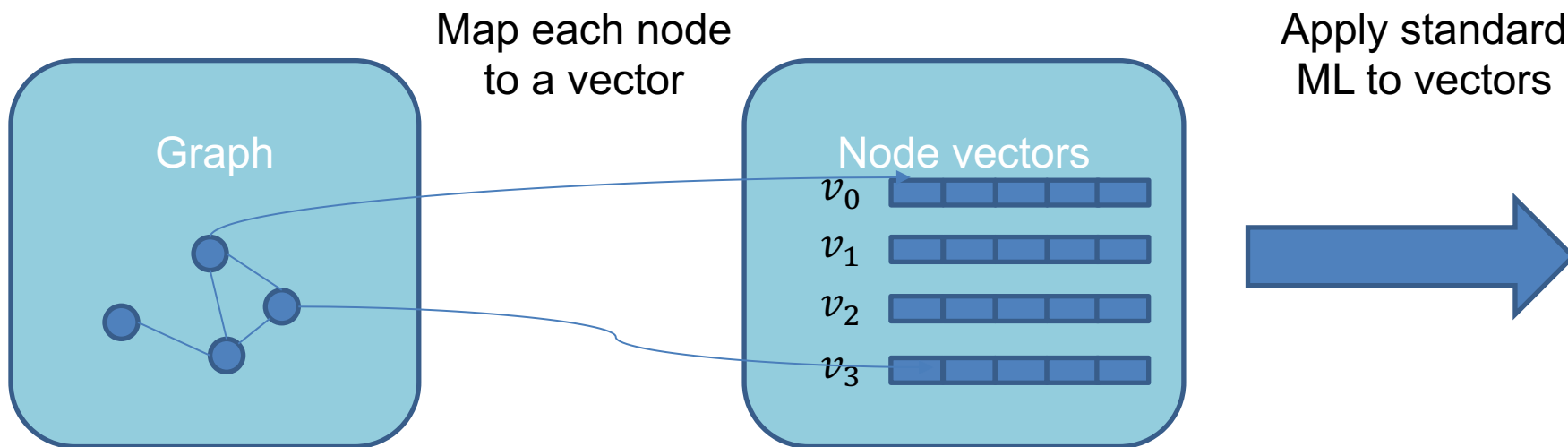Yu-Hang Tang, Oguz Selvitopi, Doru Popovici, and Aydin Buluç. A high-throughput solver for marginalized graph kernels on GPU. In Proceedings of the IPDPS, 2020.

# Graph Neural Networks

**Can be used to classify unlabeled nodes in graph**

**Want to map nodes to feature vectors**

- **Embed properties of graph, e.g. connectivity, in vectors per node**



Map each node to a vector

Apply standard ML to vectors

Graph

Node vectors

$v_0$
$v_1$
$v_2$
$v_3$

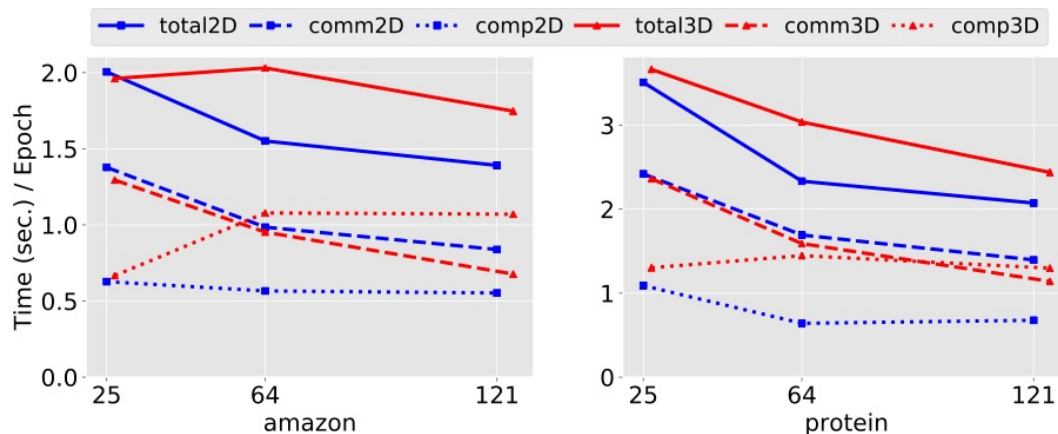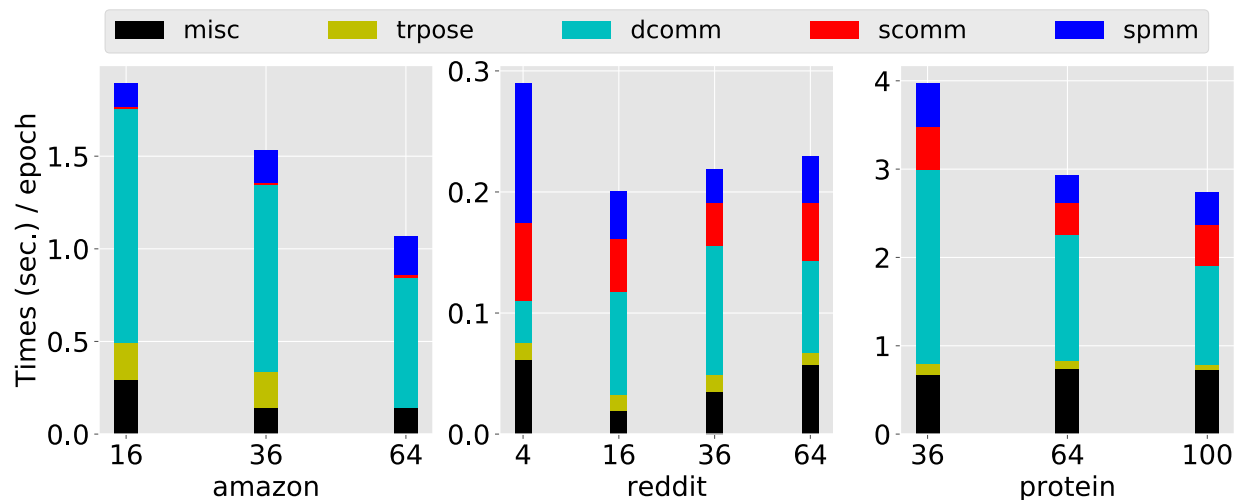How do we compute these vectors? With a GNN

# GNN Training

- Each node is initialized with a feature vector
  - $H^0$ has initial feature vector per node $(n \; x \; f)$
- Each node aggregates vectors of its neighbors, applies a weight
- Each layer computes gradients

```
for i = 1 … E
    for l = 1 … L
        Zl = AT * Hl-1 * Wl
        Hl = σ(Zl)

    ...
    for l = L-1 … 1
        Gl = A * Gl+1 * (Wl+1)T ⊙ σ'(Zl)
       dH/dW = (Hl-1)T * A * Gl
```

$A \in n \; x \; n$

$H^l \in n \; x \; f^l$
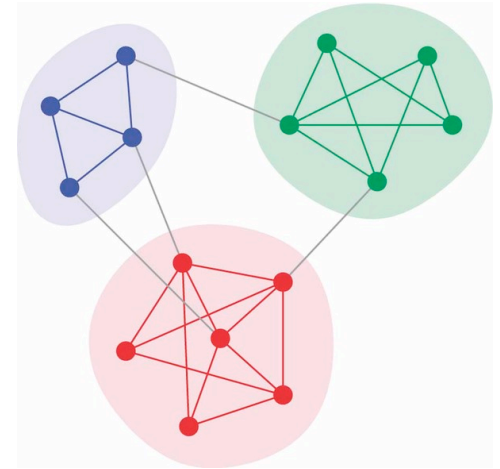
$G^l \in n \; x \; f^l$

$W^l \in f^{l-1} \; x \; f^l$

- A is sparse and f << n, so the main workhorse is SpMM (sparse matrix times tall–skinny dense matrix)

# Communication avoidance in GNN Training

Alok Tripathy, Katherine Yelick, Aydın Buluç. Reducing Communication in Graph Neural Network Training. arXiv

# The Markov Cluster Algorithm (MCL)

Widely popular and successful algorithm for discovering clusters (e.g. protein families) in protein interaction and protein sequence similarity networks

The number of **edges or higher-length paths** between two arbitrary nodes in a cluster is greater than the number of paths between nodes from different clusters

**Random walks** on the graph will frequently remains within a cluster

The algorithm **computes the probability** of random walks through the graph and **removes lower probability terms** to form clusters

# The Markov Cluster Algorithm (MCL)



**Initial network**      **Iteration 1**      **Iteration 2**      **Iteration 3**

**At each iteration:**

**Step 1** (Expansion): Squaring the matrix while
        pruning (a) small entries, (b) denser columns
**Naïve implementation:** sparse matrix-matrix product (SpGEMM),
followed by column-wise top-K selection and column-wise pruning
**Step 2** (Inflation) : taking powers entry-wise

# A combined expansion and pruning step



- ❑ b: number of columns in the output constructed at once
  - – Smaller b: less parallelism, memory efficient (b=1 is equivalent to sparse matrix-sparse vector multiplication used in MCL)
  - – Larger b: more parallelism, memory intensive

# HipMCL: High-performance MCL

- MCL process is both **computationally expensive** and **memory hungry**, limiting the sizes of networks that can be clustered
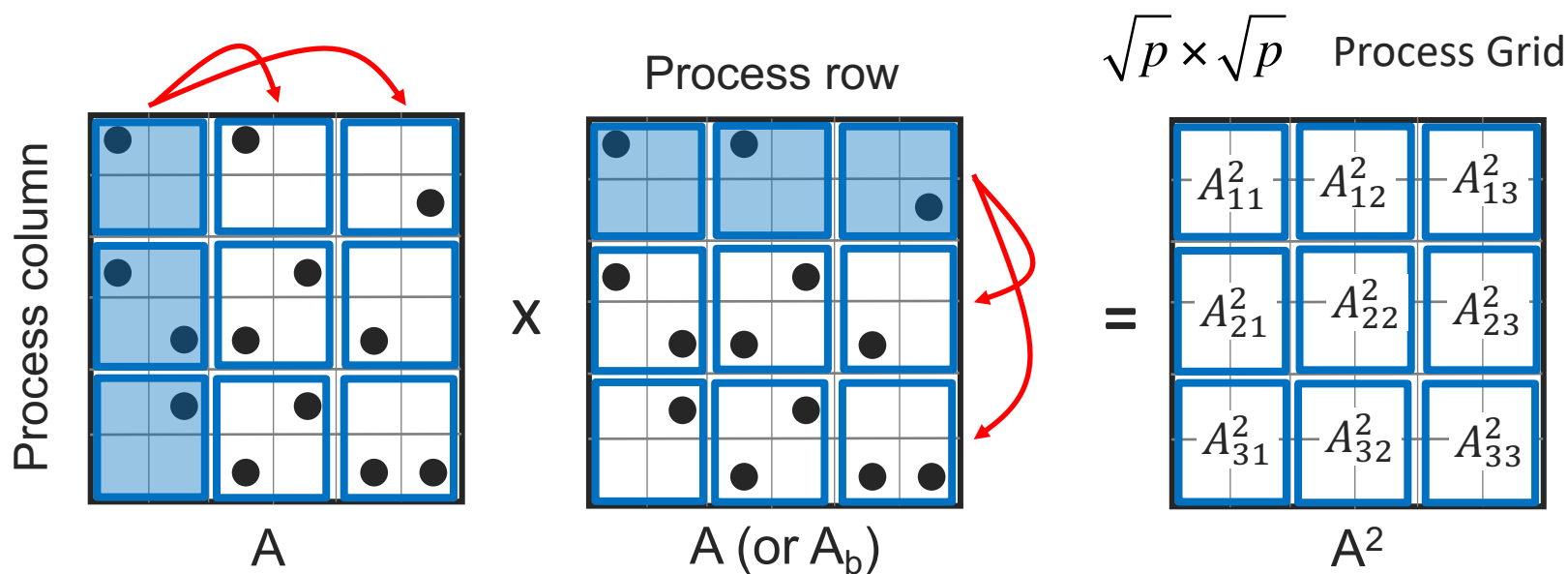- HipMCL overcomes such limitation via **sparse parallel algorithms**.
- **Up to 1000X times faster** than original MCL with same accuracy.

$$\sqrt{p} \times \sqrt{p} \quad \text{Process Grid}$$

| | | |
|---|---|---|
| $A^2_{11}$ | $A^2_{12}$ | $A^2_{13}$ |
| $A^2_{21}$ | $A^2_{22}$ | $A^2_{23}$ |
| $A^2_{31}$ | $A^2_{32}$ | $A^2_{33}$ |

Process row

Process column

A          X          A (or $A_b$)          =          $A^2$

A. Azad, G. Pavlopoulos, C. Ouzounis, N. Kyrpides, A. Buluç; HipMCL: a high-performance parallel implementation of the Markov clustering algorithm for large-scale networks, *Nucleic Acids Research, 2018*

# HipMCL on large networks

| Data | Proteins | Edges | #Clusters | HipMCL time | platform |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Isolate-1 | 47M | 7 B | 1.6M | 1 hr | 1024 nodes Edison |
| Isolate-2 | 69M | 12 B | 3.4M | 1.66 hr | 1024 nodes Edison |
| Isolate-3 | 70M | 68 B | 2.9M | 2.41 hr | 2048 nodes Cori KNL |
| MetaClust50 | 282M | 37B | 41.5M | 3.23 hr | 2048 nodes Cori KNL |

## MCL can not cluster these networks

# HipMCL on Supercomputers with accelerators

- Recent top supercomputers are all accelerated (e.g. with GPUs)
- This is what a ORNL Summit node looks like
- There are 4608 such nodes in the system
- Challenges: (1) Utilizing all GPUs, (2) hiding the communication
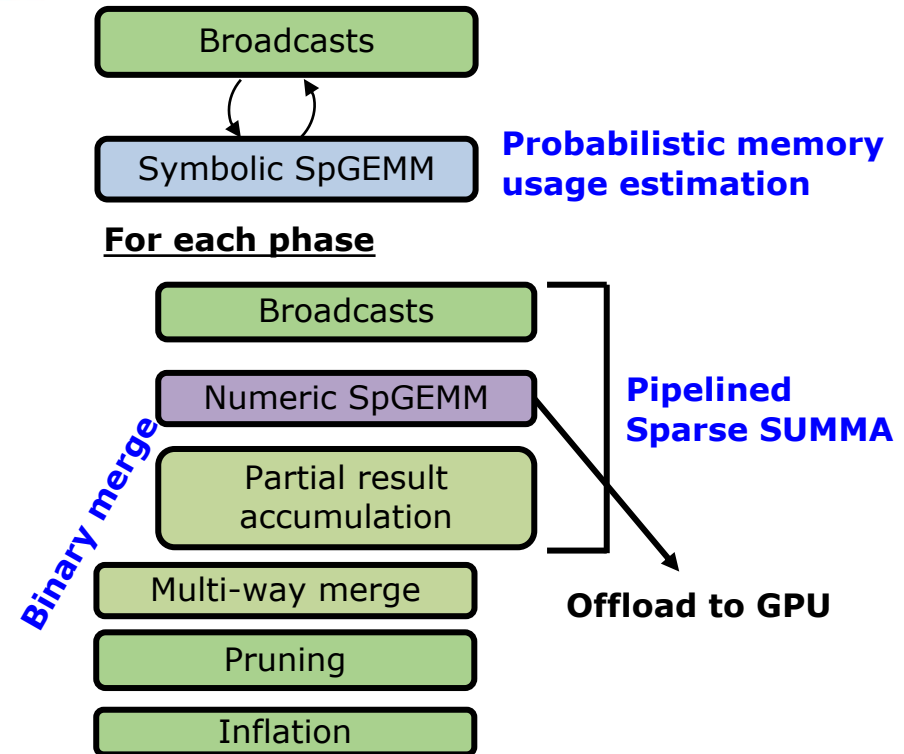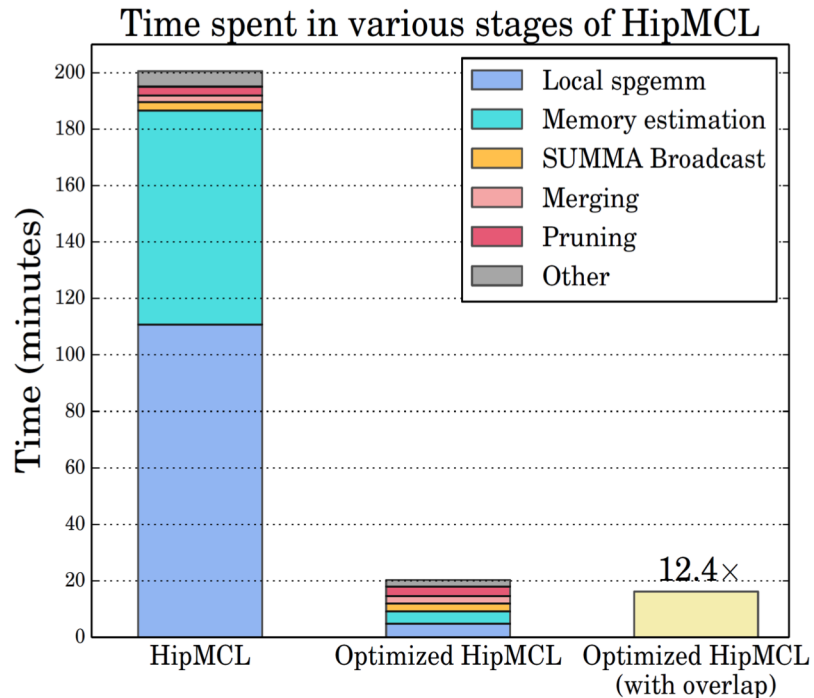


**Pipelined Sparse SUMMA**

Joint CPU-GPU distributed memory expansion of MCL algorithm

# HipMCL on Supercomputers with accelerators



Time spent in various stages of HipMCL

Legend:
- Local spgemm
- Memory estimation
- SUMMA Broadcast
- Merging
- Pruning
- Other

12.4×

Workflow diagram:
- Broadcasts
- Symbolic SpGEMM — **Probabilistic memory usage estimation**

**For each phase**

- Broadcasts
- Numeric SpGEMM
- Partial result accumulation
- Multi-way merge
- Pruning
- Inflation

**Binary merge**

**Pipelined Sparse SUMMA**

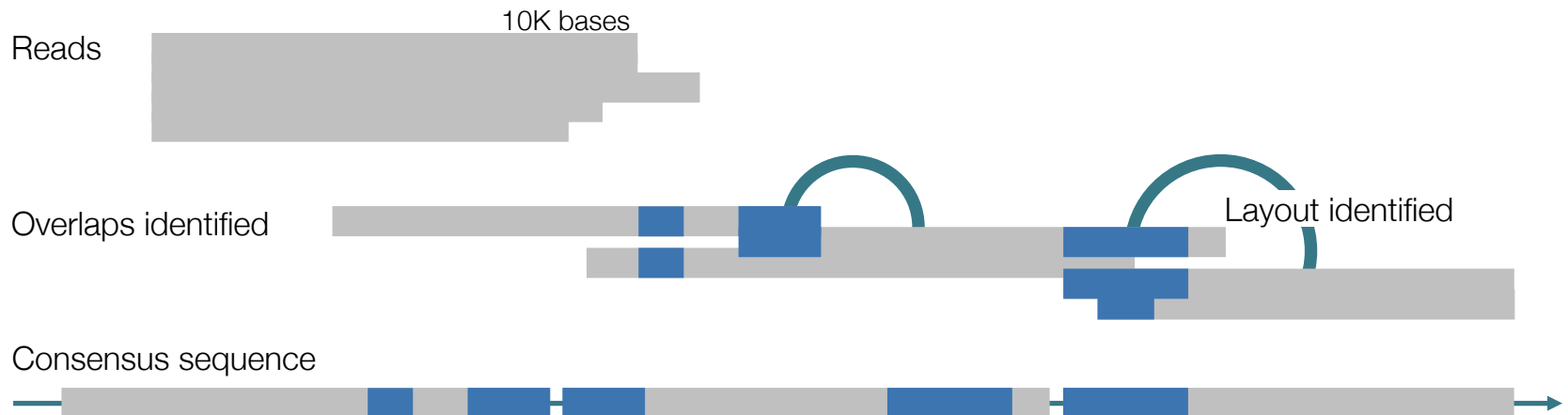**Offload to GPU**

## Other changes to HipMCL for the CPU-GPU workflow:

- *Randomized memory estimation algorithm* avoids symbolic phase
- New *eager binary merging* reduces memory footprint
- Integration of a much faster hash-based CPU SpGEMM algorithm

O. Selvitopi, M.T. Hussain, A. Azad, and A. Buluç. Optimizing high performance Markov clustering for pre-exascale architectures. IPDPS, 2020
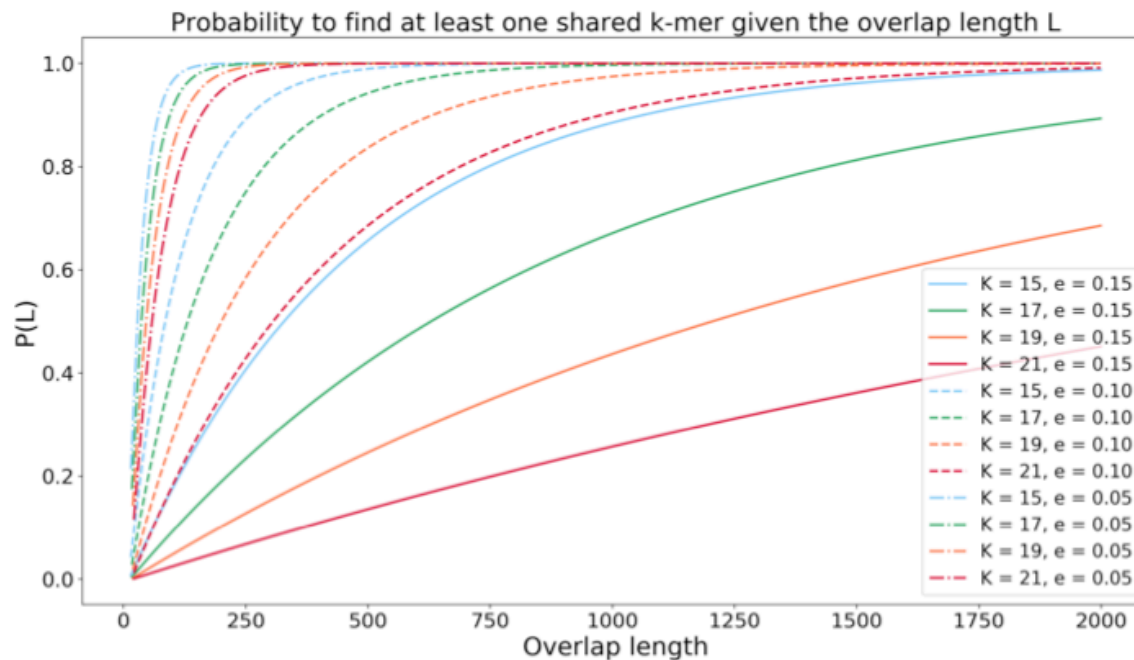
# SpGEMM for DNA read overlapping

- **Long reads** from PacBio and Oxford Nanopore have the potential to revolutionize de-novo assembly
- **Overlap-Consensus-Layout** paradigm is more suitable than de Bruijn graph paradigm.
- **Overlapping** is the most computationally expensive step.

Overlap-Layout-Consensus

Reads
10K bases

Overlaps identified
Layout identified

Consensus sequence

# SpGEMM for DNA read overlapping

- We need to quickly determine pairs of reads that are *likely to* overlap, without resorting to $O(n^2)$ comparisons

- If two reads do not share any subsequence of length k (aka a k-mer) for a reasonably short k, then they are unlikely to overlap

Probability to find at least one shared k-mer given the overlap length L

# SpGEMM for DNA read overlapping
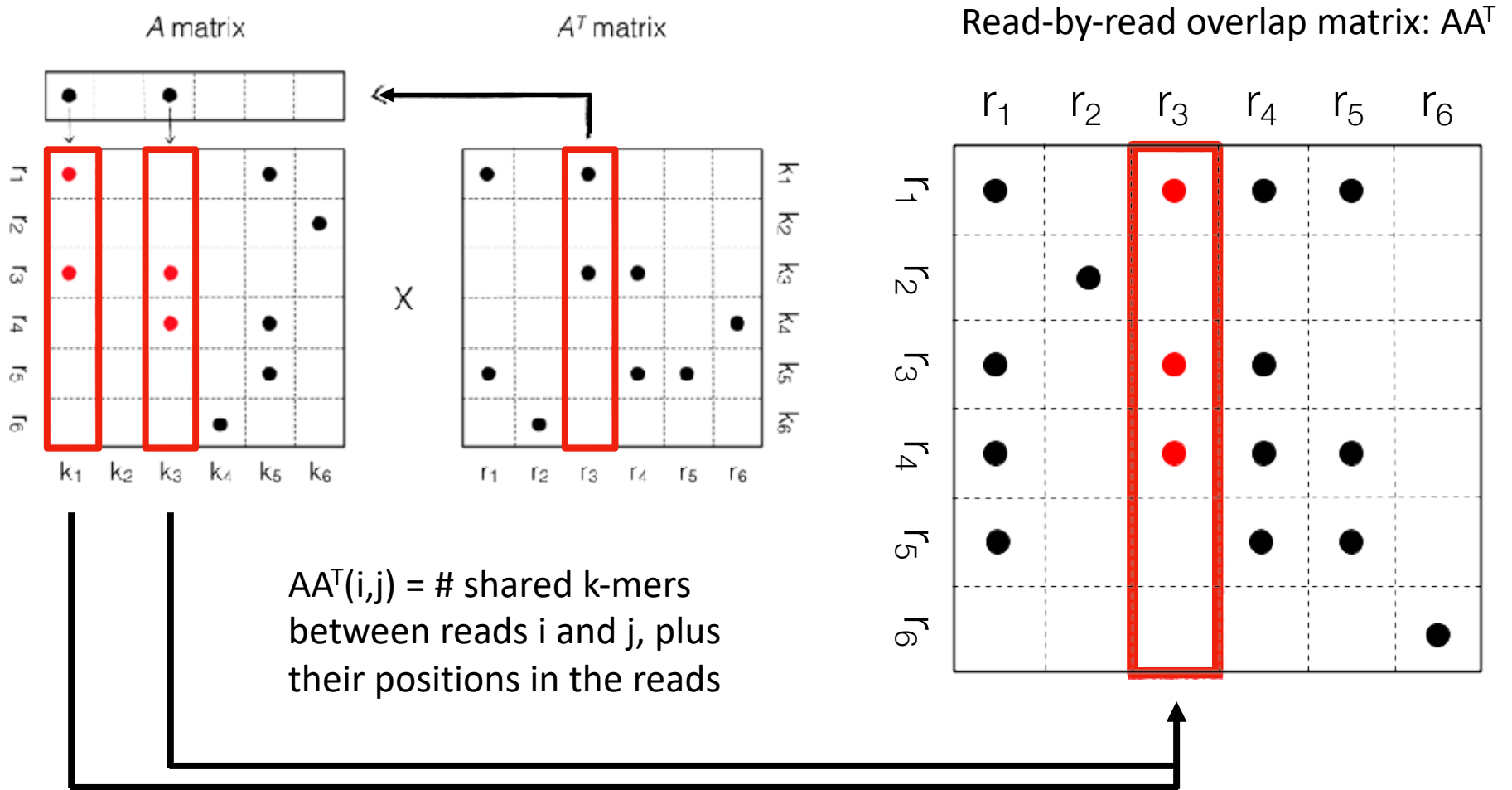
## *A* matrix



- Suppose you have counted k-mers and only retained *reliable* k-mers
- Now you can generate this ***read-by-kmer*** sparse matrix **A**
- These are all linear time computations so far

$r_i$ = $i^{th}$ read

$k_j$ = $j^{th}$ reliable *k*-mer

A(i,j) = presence of $j^{th}$ reliable *k*-mer in $i^{th}$ read, plus its position

Giulia Guidi, Marquita Ellis, Daniel Rokhsar, Kathy Yelick, Aydın Buluç, BELLA: Berkeley Efficient Long-read to Long-Read Overlapper and Aligner, Biorxiv, 2018

# SpGEMM for DNA read overlapping



A matrix

$A^T$ matrix

X

Read-by-read overlap matrix: $AA^T$

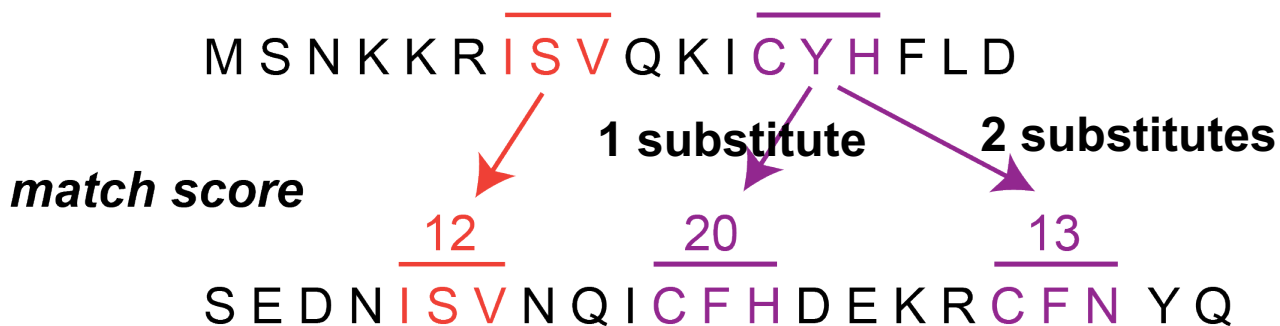$AA^T(i,j)$ = # shared k-mers between reads i and j, plus their positions in the reads

Use any fast SpGEMM algorithm, as long as it runs on *arbitrary semirings*

# SpGEMM for many-to-many protein alignment

- Idea similar to BELLA, but removing the exact match restriction
- For homology detection, need to catch weaker signal (~30% ANI)
- K-mers with substitutes may be more valuable than exact matches!



BLOSUM 62 scoring matrix

(positive values are shaded)



M S N K K R I S V Q K I C Y H F L D

1 substitute          2 substitutes

match score

12          20          13

S E D N I S V N Q I C F H D E K R C F N Y Q

# SpGEMM for many-to-many protein alignment

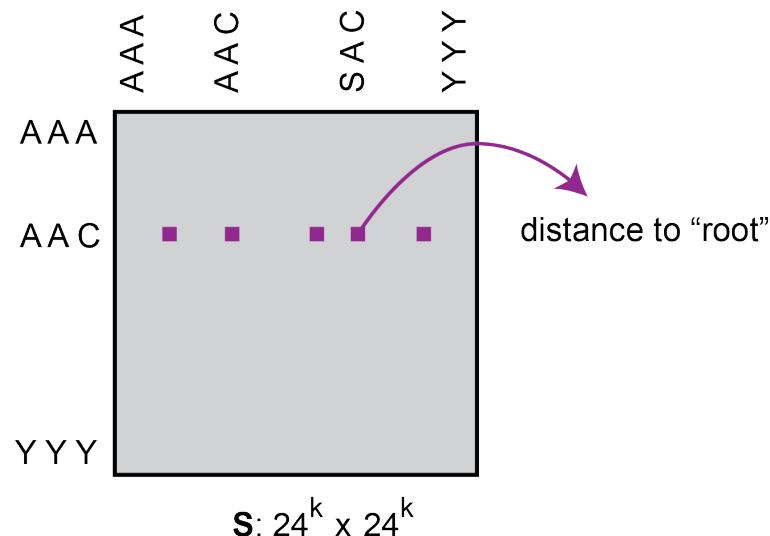## Introduce new sparse matrix S

- **Contains substitution information**

- **Each entry**

    - **Substitution cost**

**Exact k-mers $\rightarrow$ C=AA$^T$**

**Substitute k-mers $\rightarrow$ C=ASA$^T$**

**New semiring**



distance to "root"

**S**: $24^k$ x $24^k$

Oguz Selvitopi, Saliya Ekanayake, Giulia Guidi, Georgios Pavlopoulos, Ariful Azad, and Aydın Buluç. Distributed Many-to-Many Protein Sequence Alignment Using Sparse Matrices. SC'20.

# Acknowledgments

Ariful Azad, Tim Davis, Saliya Ekanayake, Marquita Ellis, John Gilbert, Giulia Guidi, Jeremy Kepner, Nikos Krypides, Tim Mattson, Scott McMillan, Jose Moreira, John Owens, Georgios Pavlopoulos, Dan Rokhsar, Oguz Selvitopi, Yu-Hang Tang, Alok Tripathy, Carl Yang, Kathy Yelick.

- My Research Team: http://passion.lbl.gov
- Our (new) Youtube Channel: http://shorturl.at/lpFRY
- The GraphBLAS Forum: http://graphblas.org